

# PATENT ABSTRACTS OF JAPAN

(11)Publication number : 11-053400

(43)Date of publication of application : 26.02.1999

BEST AVAILABLE COPY

(51)Int.Cl. G06F 17/30  
G06F 12/00

(21)Application number : 09-220233

(71)Applicant : NEC CORP

(22)Date of filing : 31.07.1997

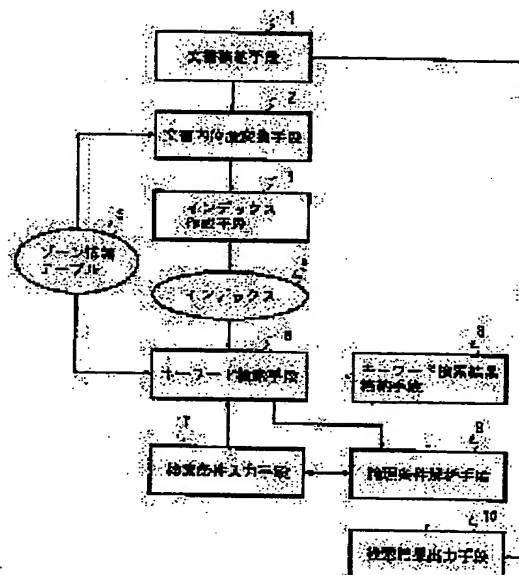
(72)Inventor : AKAMINE SUSUMU

## (54) STRUCTURED DOCUMENT RETRIEVAL DEVICE AND MACHINE READABLE RECORDING MEDIUM FOR RECORDING PROGRAM

### (57)Abstract:

**PROBLEM TO BE SOLVED:** To execute zone retrieval at high speed in a structured document retrieval device.

**SOLUTION:** A zone information table 3 commonly defines and holds a zone name for each zone and a range of a position that the zone can occupy in the whole retrieval object document. A conversion means 2 of a position in document refers to the zone information table 3 at the time of preparing an index and prepares a zone position conversion document that moves a character string of each zone of an original document to a position indicated by the zone information table 3. Thus, it becomes possible to identify from within the document in which zone the character string exists. An index preparation means 4 makes the zone position conversion document prepared by the conversion means 2 of a position in document and prepares an index 5. A keyword retrieval means 6 executes zone retrieval only by selecting the one in which an appearance position of the keyword corresponds to a zone name of a zone instructed to be a retrieval object by a user at the time of retrieving the keyword.



### LEGAL STATUS

[Date of request for examination] 31.07.1997

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number] 2962287

[Date of registration] 06.08.1999

[Number of appeal against examiner's decision of rejection]

(51) IntCl.<sup>6</sup>

識別記号

F I

G 0 6 F 17/30  
12/00

5 4 7

G 0 6 F 15/419  
12/00  
15/403 2 0  
5 4 7 H  
3 7 0 A

審査請求 有 請求項の数 5 F D (全 11 頁)

(21) 出願番号

特願平9-220233

(22) 出願日

平成9年(1997) 7月31日

(71) 出願人 000004237

日本電気株式会社

東京都港区芝五丁目7番1号

(72) 発明者 赤峯 亨

東京都港区芝五丁目7番1号 日本電気株式会社内

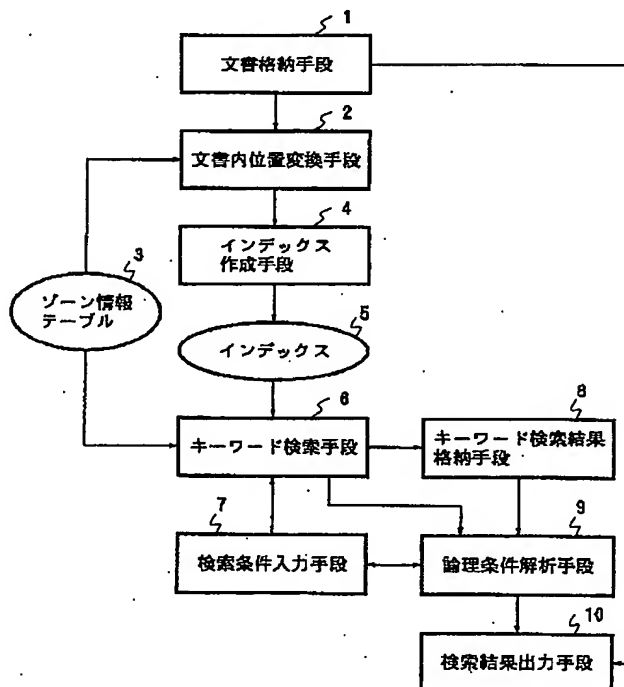
(74) 代理人 弁理士 境 廣巳

(54) 【発明の名称】 構造化文書検索装置及びプログラムを記録した機械読み取り可能な記録媒体

(57) 【要約】

【課題】 構造化文書検索装置に於いて、ゾーン検索を高速に行う。

【解決手段】 ゾーン情報テーブル3には、各ゾーンのゾーン名と、そのゾーンが取り得る位置の範囲を全検索対象文書で共通に定めて保持している。文書内位置変換手段2は、インデックス作成時に、ゾーン情報テーブル3を参照して、オリジナルの文書の各ゾーンの文字列を、ゾーン情報テーブル3によって示される位置に移動したゾーン位置変換文書を作成する。これにより、文書内位置から文字列がどのゾーンに存在するかの識別が可能になる。インデックス作成手段4は、文書内位置変換手段2により作成されたゾーン位置変換文書を対象にしてインデックス5を作成する。キーワード検索手段6は、キーワードを検索する際に、キーワードの出現位置がユーザによって検索対象とすることが指示されたゾーンのゾーン名に対応するものだけを選択することで、ゾーン検索を行う。



## 【特許請求の範囲】

【請求項 1】 複数のゾーンから構成される構造化文書が複数格納された文書格納手段と、  
ゾーン位置変換文書に於ける各ゾーンの位置を示す情報が格納されたゾーン情報テーブルと、  
前記文書格納手段に格納されている構造化文書中の各ゾーンを前記ゾーン情報テーブルの内容によって示される位置に移動させたゾーン位置変換文書を作成する文書内位置変換手段と、  
該文書内位置変換手段によって作成されたゾーン位置変換文書に基づいて、キー文字列と、そのキー文字列が存在する構造化文書の文書識別子と、そのキー文字列のゾーン位置変換文書に於ける文書内位置とが対応して格納されたインデックスを作成するインデックス作成手段と、  
検索対象にするゾーンのゾーン名とキーワードとを含む検索条件式を受け付ける検索条件入力手段と、  
該検索条件入力手段が受け付けた検索条件式中のキーワードをキーにして前記インデックスを検索し、その結果得られた前記キーワードが存在する構造化文書の文書識別子、文書内位置と、前記ゾーン情報テーブルの内容とに基づいて、前記検索条件式中のゾーン名によって示されるゾーンに前記キーワードが存在する構造化文書の文書識別子を求めるキーワード検索手段とを備えたことを特徴とする構造化文書検索装置。

【請求項 2】 前記ゾーン情報テーブルには、ゾーン名と、そのゾーン名のゾーンをゾーン位置変換文書内のどの位置に配置するのかを示すゾーン位置情報とが対応して格納されていることを特徴とする請求項 1 記載の構造化文書検索装置。

【請求項 3】 前記検索条件式は、検索対象にするゾーンのゾーン名とキーワードとから構成される検索項目が論理演算記号によって複数結合された形式を有し、  
前記キーワード検索手段は、前記検索条件入力手段が受け付けた検索条件式の各検索項目それぞれについて、その検索項目中のキーワードをキーにして前記インデックスを検索し、その検索結果と、前記ゾーン情報テーブルの内容とに基づいて、前記検索項目中のゾーン名によって示されるゾーンに前記キーワードが存在する構造化文書の文書識別子を求める構成を有し、且つ、  
前記キーワード検索手段が求めた各検索項目毎の文書識別子と、前記検索条件入力手段が受け付けた検索条件式中の各検索項目を結合する論理演算記号とに基づいて、前記検索条件式を満足させる構造化文書の文書識別子を求める論理条件解析手段を備えたことを特徴とする請求項 2 記載の構造化文書検索装置。

【請求項 4】 前記論理条件解析手段が求めた文書識別子の構造化文書を前記文書格納手段から読み出して出力する検索結果出力手段を備えたことを特徴とする請求項 3 記載の構造化文書検索装置。

【請求項 5】 複数のゾーンから構成される構造化文書が複数格納された文書格納手段と、ゾーン位置変換文書に於ける各ゾーンの位置を示すゾーン位置情報が格納されたゾーン情報テーブルとを備えたコンピュータを、  
前記文書格納手段に格納されている構造化文書中の各ゾーンを前記ゾーン情報テーブルの内容によって示される位置に移動させたゾーン位置変換文書を作成する文書内位置変換手段、  
該文書内位置変換手段によって作成されたゾーン位置変換文書に基づいて、キー文字列と、そのキー文字列が存在する構造化文書の文書識別子と、そのキー文字列の文書内位置とが対応して格納されたインデックスを作成するインデックス作成手段、  
検索対象にするゾーンのゾーン名とキーワードとを含む検索条件式を受け付ける検索条件入力手段、  
該検索条件入力手段が受け付けた検索条件式中のキーワードをキーにして前記インデックスを検索し、その結果得られた前記キーワードを含む文書の文書識別子、文書内位置と、前記ゾーン情報テーブルの内容とに基づいて、前記検索条件式中のゾーン名によって示されるゾーンに前記キーワードが存在する構造化文書の文書識別子を求めるキーワード検索手段として機能させるためのプログラムを記録した機械読み取り可能な記録媒体。

## 【発明の詳細な説明】

## 【0001】

【発明の属する技術分野】 本発明は、1 件の文書が複数の論理構造（ゾーン）から構成されている構造化文書を対象にした検索技術に関し、特に、ユーザによって指定されたゾーンのみを検索対象にして構造化文書の検索を行う技術に関する。

## 【0002】

【従来の技術】 近年、SGML (Standard Generalized Markup Language) に代表される、1 件の文書を複数のゾーンで構成した構造化文書が頻繁に用いられるようになってきている。それにつれて、構造化文書の特定ゾーンのみを検索対象にした検索（ゾーン検索）を行うことが、精度の高い検索を行う上で重要な機能になってきている。

【0003】 ゾーン検索を行う従来の技術としては、例えば、検索条件を「検索対象とするゾーンの開始タグと終了タグの間にキーワードを含む文書」とし、テキスト全体を対象にして文字列照合を行うことにより、検索条件を満足させる文書を検索するようにしたものがある。しかし、この方法は、テキスト全体を対象にして文字列照合を行うので、検索時間が非常に長くなってしまふ。このような問題点を解決するため、テキストの不要部分（検索対象とすることが指定されたゾーン以外のゾーン）をスキップして検索を行うようにした技術も提案されている（例えば、特開平 8-16600 号公報）。この技術によれば、文字列照合を行う範囲が少なくなるた

め、先の従来例に比較して検索時間を短くすることができる。しかし、ギガバイトクラスの大規模データに対する検索では、不要部分をスキップすることにより文字列照合範囲を例えば 1 0 分の 1 程度に縮小できたとしても、百メガバイトクラスのテキストを対象にして文字列照合を行うことが必要になるため、高速な検索は望めない。

【0004】このように、文字列照合によりゾーン検索を行う方法では、高速な検索を行うことが難しいため、大規模データに対する高速な検索を可能にするために作成された全文インデックスを利用してゾーン検索を行うことが考えられる。全文インデックスは、キー文字列がキー情報として格納されたキー情報部と、キー情報部に格納された各キー文字列それぞれについてそのキー文字列が存在する文書の文書識別子、文書内位置が位置情報として格納された位置情報部とから構成されるものであり、このような全文インデックスを利用してゾーン検索を行う方法としては、下記 (A) ~ (C) の 3 つの方法が考えられる。

#### 【0005】

(A) ゾーンに関する情報をキー情報部に持たせる方法。

(B) ゾーンに関する情報を位置情報部に持たせる方法。

(C) 全文インデックスとは別のゾーンに関するインデックスを作成する方法。

【0006】(A)の方法では、全文インデックスのキー情報部に、キー文字列とそのキー文字列が存在するゾーンのゾーン名とのペアからなるキー情報を格納しておく。その際、複数のゾーンに存在するキー文字列については、各ゾーン毎にキー文字列とゾーン名との対からなるキー情報を格納する。位置情報部には、各キーワード毎に該当する文書の文書識別子、文書内位置が格納される。そして、検索時には、ユーザによって指定されたゾーン名とキーワードとをキーにして全文インデックスを検索することにより、上記ゾーン名のゾーンに、上記キーワードを含む文書を探し出すようにしている。

【0007】(B)の方法では、位置情報部に格納する位置情報に、位置情報として文書識別子、文書内位置の他にゾーン名も併せ持たせておく。そして、検索時には、まず、ユーザによって指定されたキーワードをキーにして全文インデックスを検索することにより、上記キーワードを含む文書の位置情報を全て求め、その後、上記位置情報の中からユーザによって指定されたゾーン名を含む位置情報を選択することにより、ゾーン検索を行うようにしている。

【0008】(C)の方法では、全文インデックスとは別に、検索対象とする全ての文書それぞれの、各ゾーンの開始位置、終了位置が格納されたゾーン用インデックスを作成しておく。そして、検索時には、まず、全文イ

ンデックスを検索することにより、ユーザが指定したキーワードを含む文書の位置情報を取得する。その後、ゾーン用インデックスを検索し、上記文書の、ユーザによって指定されたゾーンの開始位置、終了位置を取得する。更に、位置情報中の文書内位置と取得した上記ゾーンの開始位置、終了位置とに基づいて、上記文書のユーザによって指定されたゾーン内に、ユーザによって指定されたキーワードが存在するか否かをチェックすることにより、ゾーン検索を行う (例えば、特開平 8 - 3 1 4 9 6 6 号公報)。

#### 【0009】

【発明が解決しようとする課題】しかしながら、上述した (A) の方法は、ゾーンの数に比例してキー情報数が多くなるため、全文インデックスのサイズが大きくなるという問題がある。更に、全文インデックスのサイズが大きくなることにより、検索速度が低下するという問題もある。

【0010】また、(B)の方法は、全文インデックスの位置情報部の容量が多くなるため、全文インデックスのサイズが大きくなるという問題がある。更に、位置情報部からの読み出し量が多くなるため、検索速度が低下するという問題もある。

【0011】また、(C)の方法は、検索対象とする全ての文書それぞれの、各ゾーンの開始位置、終了位置を格納したゾーン用インデックスが必要になるため、インデックスサイズが大きくなるという問題がある。更に、全文インデックスを検索することにより探し出した各該当文書について、ゾーン用インデックスを検索し、ユーザが指定したキーワードが、ユーザが指定したゾーン内に存在するか否かをチェックする必要があるため、この処理がオーバーヘッドになって検索速度が低下してしまうという問題が生じる。

【0012】そこで、本発明の目的は、全文インデックスを利用したゾーン検索に於いて、インデックスサイズを小さくし、且つ検索速度を高速化することにある。

#### 【0013】

【課題を解決するための手段】本発明の構造化文書検索装置は、上記目的を達成するため、複数のゾーンから構成される構造化文書が複数格納された文書格納手段と、ゾーン位置変換文書に於ける各ゾーンの位置を示す情報が格納されたゾーン情報テーブルと、前記文書格納手段に格納されている構造化文書中の各ゾーンを前記ゾーン情報テーブルの内容によって示される位置に移動させたゾーン位置変換文書を作成する文書内位置変換手段と、該文書内位置変換手段によって作成されたゾーン位置変換文書に基づいて、キー文字列と、そのキー文字列が存在する構造化文書の文書識別子と、そのキー文字列のゾーン位置変換文書に於ける文書内位置とが対応して格納されたインデックスを作成するインデックス作成手段と、検索対象にするゾーンのゾーン名とキーワードとを

含む検索条件式を受け付ける検索条件入力手段と、該検索条件入力手段が受け付けた検索条件式中のキーワードをキーにして前記インデックスを検索し、その結果得られた前記キーワードが存在する構造化文書の文書識別子、文書内位置と前記ゾーン情報テーブルの内容とに基づいて、前記検索条件式中のゾーン名によって示されるゾーンに前記キーワードが存在する構造化文書の文書識別子を求めるキーワード検索手段とを備えたものである。

【0014】この構成に於いては、ゾーン検索を行うための準備として、文書内位置変換手段が、文書格納手段に格納されている構造化文書中の各ゾーンをゾーン情報テーブルの内容によって示される位置に移動させたゾーン位置変換文書を作成し、インデックス作成手段が、文書内位置変換手段によって作成されたゾーン位置変換文書に基づいて、キー文字列と、そのキー文字列が存在する構造化文書の文書識別子と、そのキー文字列のゾーン位置変換文書に於ける文書内位置とが対応して格納されたインデックスを作成する。

【0015】そして、ゾーン検索時に、ユーザが検索対象にするゾーンのゾーン名とキーワードとを含む検索条件式を入力すると、検索条件入力手段がそれを受け付け、キーワード検索手段が上記検索条件式中のキーワードをキーにしてインデックスを検索し、その結果得られた前記キーワードが存在する構造化文書の文書識別子、文書内位置とゾーン情報テーブルの内容とに基づいて、上記検索条件式中のゾーン名によって示されるゾーンに前記キーワードが存在する構造化文書の文書識別子を求める。

【0016】

【発明の実施の形態】次に本発明の実施の形態について図面を参照して詳細に説明する。

【0017】図1は本発明の実施例のブロック図であり、文書格納手段1と、文書内位置変換手段2と、ゾーン情報テーブル3と、インデックス作成手段4と、インデックス5と、キーワード検索手段6と、検索条件入力手段7と、キーワード検索結果格納手段8と、論理条件解析手段9と、検索結果出力手段10とから構成されている。

【0018】文書格納手段1には、検索対象となる複数の構造化文書が格納されている。図2は文書格納手段1の内容例を示した図である。本実施例では、説明を簡単にするため、文書格納手段1には、文書識別子ID1、ID2の2つの構造化文書21、22が格納されているとする。各構造化文書21、22は、それぞれゾーン名「全体」、「発明の名称」、「要約」、「目的」、「構成」の各ゾーンから構成され、各ゾーンは、ゾーン開始タグ、ゾーン終了タグ等により分割されている。

【0019】ゾーン情報テーブル3には、文書内位置変換手段2で作成するゾーン位置変換文書に於ける各ゾー

ンの位置を示す情報が格納されている。図3はゾーン情報テーブル3の内容例を示した図であり、ゾーン名と、そのゾーン名のゾーンをゾーン位置変換文書内のどの位置に配置するのかを示すゾーン位置情報とが格納されている。図3の例は、ゾーン名「全体」、「発明の名称」、「要約」、「目的」、「構成」のゾーンを、それぞれゾーン位置変換文書内の「第1文字目～第2000文字目」、「第1文字目～第500文字目」、「第501文字目～第2000文字目」、「第501文字目～第1000文字目」、「第1001文字目～第2000文字目」に配置することを示している。

【0020】文書内位置変換手段2は、ゾーン情報テーブル3を参照し、文書格納手段1に格納されている構造化文書に対して、各ゾーンの文字列をゾーン情報テーブル3中のゾーン位置情報によって示される位置に移動したゾーン位置変換文書を作成する機能を有する。従って、各構造化文書では異なる位置に存在していた各ゾーンの文字列は、文書内位置変換手段2によって作成されたゾーン位置変換文書では、同じ範囲に存在することになる。

【0021】インデックス作成手段4は、文書内位置変換手段2で作成された各構造化文書に対応するゾーン位置変換文書に基づいてインデックス5を作成する機能を有する。インデックス5は、図4に示すように、キー情報部51と、位置情報部52とを有している。キー情報部51にはN文字組や単語等のキー情報が格納され、位置情報部51には、キー情報が存在する構造化文書の文書識別子と、そのキー情報が存在するゾーン位置変換文書内の位置とが格納される。ここで、位置情報「1-j」は、文書識別子ID1の構造化文書に対応するゾーン位置変換文書の第j文字目を表している。従って、図4の例の第1番目のエリアは、文字「文」が文書識別子ID1の構造化文書21中に存在し、それと対応するゾーン位置変換文書21'では第1文字目、第504文字目に存在することと、文字「文」が文書識別子ID2の構造化文書22中に存在し、それと対応するゾーン位置変換文書22'では第1文字目、第501文字目に存在することを表している。

【0022】検索条件入力手段7は、ユーザによって入力された検索条件式を受け付ける機能、論理条件解析手段9を利用して検索条件式を検索項目に分解する機能、検索項目をキーワード検索手段6に渡す機能等を有する。ユーザが入力する検索条件式は、検索対象とするゾーン名とキーワードとのペアからなる検索項目を1つ或いは複数含むものであり、検索項目を複数含む場合は、各検索項目は、AND、OR等の論理演算記号によって結合される。図5は、ユーザが入力する検索条件式の1例を示した図であり、2つの検索項目が論理演算記号ANDによって結合されている。この検索条件式は、ゾーン名「発明の名称」のゾーンにキーワード「検索」を含

み、且つゾーン名「要約」のゾーンにキーワード「インデックス」を含む構造化文書の検索を指示するものである。

【0023】キーワード検索手段6は、検索条件入力手段7から渡された各検索項目中のキーワードをキーにしてインデックス5を検索することにより、各検索項目それぞれについて、その検索項目中のキーワードが現れる文書の文書識別子、文書内位置を全て求める機能や、各検索項目それぞれについて、ゾーン情報テーブル3を参照して検索項目に含まれているゾーン名によって示されるゾーンのゾーン位置を求める機能や、各検索項目それぞれについて、その検索結果（文書識別子、文書内位置）の中に上記ゾーン位置内の位置を示す検索結果があれば、その検索結果中の文書識別子とそれが何番目の検索項目についてのものなのかを示す情報とをペアにしてキーワード検索結果格納手段8に格納する機能等を有する。

【0024】論理条件解析手段9は、キーワード検索結果格納手段8に格納されている検索項目毎の検索結果

（文書識別子）と、検索条件入力手段7が受け付けた検索条件式中の各検索項目を結合する論理演算記号とに基づいて、上記検索条件式を満足させる構造化文書の文書識別子を求める機能を有する。

【0025】検索結果出力手段10は、論理条件解析手段9が求めた文書識別子を有する構造化文書を文書格納手段1から取り出し、プリンタ、CRT等の出力装置（図示せず）に出力する機能を有する。

【0026】図6は文書内位置変換手段2の処理例を示す流れ図、図7はインデックス作成手段4の処理例を示す流れ図、図8は検索条件入力手段7の処理例を示す流れ図、図9は検索条件入力手段7から検索条件式が渡されたときの論理条件解析手段9の処理例を示す流れ図、図10はキーワード検索手段6から終了通知が送られてきたときの論理条件解析手段9の処理例を示す流れ図、図11はキーワード検索手段6の処理例を示す流れ図であり、以下各図を参照して本実施例の動作を説明する。

【0027】先ず、インデックス5の作成時の動作を説明する。

【0028】インデックス5の作成時、文書内位置変換手段2は、図6の流れ図に示すように、文書格納手段1から未処理の構造化文書を1つ入力する（S61）。その後、文書内位置変換手段2は、入力した構造化文書の各ゾーン中の文字列をゾーン情報テーブル3中のゾーン位置情報によって示される位置に移動させたゾーン位置変換文書を作成し（S63）、作成したゾーン位置変換文書とS61で入力した構造化文書の文書識別子とをインデックス作成手段4に渡す（S64）。以上の処理を未処理の構造化文書がなくなるまで（S62がNO）、繰り返す。

【0029】本実施例の場合、文書格納手段1には図2

に示すような文書識別子ID1、ID2の構造化文書21、22が格納され、更に、ゾーン情報テーブル3の内容は図3に示すものになっているので、文書内位置変換手段2は、図12に示すようなゾーン位置変換文書21'、22'を順次作成してインデックス作成手段4に渡すことになる。

【0030】つまり、文書内位置変換手段2は、図3に示したゾーン情報テーブル3の内容を参照し、その内容に従って、構造化文書21のゾーン「発明の名称」中に存在する文字列「文書検索装置」を第1文字目から始まる位置に移動させ、ゾーン「目的」中に存在する文字列「高速に文書を検索する。」を第501文字目から始まる位置に移動させ、ゾーン「構成」中に存在する文字列「インデックス作成手段と……。」を第1001文字目から始まる位置に移動させた図12に示すようなゾーン位置変換文書21'を作成して文書識別子ID1と共にインデックス作成手段4に渡す。同様に、文書内位置変換手段2は、構造化文書22のゾーン「発明の名称」中に存在する文字列「文書処理装置」を第1文字目から始まる位置に移動させ、ゾーン「目的」中に存在する文字列「文書を……。」を第501文字目から始まる位置に移動させ、ゾーン「構成」中に存在する文字列「検索手段と……。」を第1001文字目から始まる位置に移動させた図12に示すようなゾーン位置変換文書22'を作成して文書識別子ID2と共にインデックス作成手段4に渡す。

【0031】図12から判るように、変換処理後の各ゾーン位置変換文書21'、22'は、ゾーン名「全体」、「発明の名称」、「要約」、「目的」、「構成」の各ゾーンが、ゾーン情報テーブル3中のゾーン位置情報によって示される位置に必ず存在することになる。

【0032】インデックス作成手段4は、文書内位置変換手段2からゾーン位置変換文書、文書識別子が渡されると、図7の流れ図に示すように、ゾーン位置変換文書の先頭位置に注目する（S71）。そして注目位置に、インデックス5のキー情報部51に格納すべき文字が存在するか否かを判断する（S72）。格納すべき文字か否かの判断は、例えば、空白文字、句読点等、格納する必要のない文字を予め定めておき、注目位置に存在する文字がそれ以外の文字であるか否かを判断することにより行う。

【0033】そして、S72に於いて、格納すべき文字が注目位置に存在すると判断した場合（S72がYES）は、その文字が既に格納済みか否かを判断する（S73）。格納済みでないと判断した場合は、注目位置に存在する文字をインデックス5のキー情報部51に格納すると共に、文書内位置変換手段2から渡された文書識別子と文書内位置（現在の注目位置）とからなる位置情報を位置情報部52に格納する（S73がNO、S74）。これに対して、格納済みであると判断した場合



は、位置情報部 5 2 に文書内位置変換手段 2 から渡された文書識別子と文書内位置とからなる位置情報を位置情報部 5 2 に格納する (S 7 3 が YES, S 7 5)。

【0034】S 7 4, 7 5 の処理が終了すると、インデックス作成手段 4 は、注目位置を次の位置に移し (S 7 6)、前述したと同様の処理を行う。また、S 7 2 で格納すべき文字が注目位置に存在しないと判断した場合も、S 7 6 の処理を行う。

【0035】以上の処理を文書内位置変換手段 2 から渡されたゾーン位置変換文書の終わりまで (S 7 7 が YES)、繰り返し行う。

【0036】本実施例の場合、インデックス作成手段 4 には、図 1 2 に示すようなゾーン位置変換文書 2 1', 文書識別子 I D 1 と、ゾーン位置変換文書 2 2', 文書識別子 I D 2 とが渡されるので、インデックス作成手段 4 に於いては、次のような処理が行われることになる。

【0037】文書内位置変換手段 2 からゾーン位置変換文書 2 1' と文書識別子 I D 1 とが渡された場合は、インデックス作成手段 4 は、先頭位置に注目したときに、文字「文」をキー情報部 5 1 に格納し、位置情報「1-1」を位置情報部 5 2 に格納する (S 7 1, S 7 4)。また、インデックス作成手段 4 は、ゾーン位置変換文書 2 1' 中の次に位置 (第 2 文字目) に注目したときは、注目位置に存在する文字「書」をキー情報部 5 1 に格納し、位置情報「1-2」を位置情報部 5 2 に格納する (S 7 4)。

また、例えば、注目位置をゾーン位置変換文書 2 1' の第 5 0 4 文字目にしたときは、注目位置に存在する文字「文」は既に格納済みであるので、位置情報部 5 2 中の上記文字「文」に対応するエントリに位置情報「1-5 0 4」を格納することになる (S 7 5)。このような処理を、ゾーン位置変換文書 2 1' の終わりまで行う。ゾーン位置変換文書 2 2' と文書 I D とが渡された場合も、インデックス作成手段 4 は前述したと同様の処理を行う。この結果、インデックス 5 の内容は、図 4 に示すものとなる。

【0038】尚、ここでは、説明を簡単に行うため、キー情報部 5 1 に格納する文字列の文字長を 1 文字としたが、これに限られるものではなく、文字長が 2 以上の N 文字組でも、単語であっても構わない。

【0039】次に、ゾーン検索時の動作について説明する。

【0040】ゾーン検索を行う場合、ユーザは、検索対象とするゾーンのゾーン名とキーワードとのペアからなる検索項目を 1 つ或いは複数含む検索条件式を検索条件入力手段 7 に入力する。前述したように、検索項目を複数含む検索条件式の場合は、各検索項目は、AND, OR 等の論理演算記号によって結合されている。

【0041】今、例えば、ユーザが検索条件式として、図 5 に示した検索条件式「(発明の名称=検索) AND (要約=インデックス)」を検索条件入力手段 7 に入力

したとする。この検索条件式は、前述したように、ゾーン名「発明の名称」の部分に文字列「検索」が現れ、且つゾーン名「要約」の部分に文字列「インデックス」が現れる構造化文書の検索を指示するものである。

【0042】ユーザが検索条件式「(発明の名称=検索) AND (要約=インデックス)」を入力すると、検索条件入力手段 7 は、図 8 の流れ図に示すように、それを受け付け、論理条件解析手段 9 に渡す (S 8 1, S 8 2)。

【0043】論理条件解析手段 9 は、検索条件式「(発明の名称=検索) AND (要約=インデックス)」が渡されると、図 9 の流れ図に示すように、検索条件式を第 1 番目の検索項目「発明の名称=検索」と、第 2 番目の検索項目「要約=インデックス」との 2 つの検索項目に分割し、それらを検索条件入力手段 7 に返す (S 9 1, S 9 2)。

【0044】検索条件入力手段 7 は、論理条件解析手段 9 から第 1 番目、第 2 番目の検索項目「発明の名称=検索」、「要約=インデックス」を受け取ると、それらをキーワード検索手段 6 に渡す (図 8 の S 8 3, S 8 4)。

【0045】キーワード検索手段 6 は、検索条件入力手段 7 から第 1 番目、第 2 番目の検索項目「発明の名称=検索」、「要約=インデックス」が渡されると、図 1 1 の流れ図に示すように、その内の 1 つに注目する (S 1 1 1)。

【0046】今、例えば、第 1 番目の検索項目「発明の名称=検索」に注目したとすると、キーワード検索手段 6 は、先ず、第 1 番目の検索項目「発明の名称=検索」中のキーワード「検索」をキーにしてインデックス 5 を検索することにより、キーワード「検索」が現れるゾーン位置変換文書の文書識別子と、文書内位置とを求める (S 1 1 3)。本実施例の場合、インデックス 5 の内容は、図 5 に示すものになっているので、S 1 1 3 を行うことにより、キーワード「検索」が、文書識別子 I D 1 のゾーン位置変換文書 2 1' の第 3 文字目～第 4 文字目と、文書識別子が I D 2 のゾーン位置変換文書 2 2' の第 1 0 0 1 文字目～第 1 0 0 2 文字目に現れることが求められる。

【0047】その後、キーワード検索手段 6 は、ゾーン情報テーブル 3 を参照し、第 1 番目の検索項目中のゾーン名「発明の名称」によって示されるゾーンのゾーン位置を求める (S 1 1 4)。本実施例の場合、「第 1 文字目～第 5 0 0 文字目」がゾーン位置として求められる (図 3 参照)。

【0048】次いで、キーワード検索手段 6 は、S 1 1 3 の検索結果の中に、S 1 1 4 で求めたゾーン位置内の位置を示すものがあれば、その検索結果中の文書識別子とそれが第 1 番目の検索項目についてのものであることを示す情報とをペアにしてキーワード検索結果格納手段

8に格納する(S115)。この例では、S113の検索結果が、「文書識別子ID1のゾーン位置変換文書21'の第3文字目～第4文字目」、「文書識別子がID2のゾーン位置変換文書22'の第1001文字目～第1002文字目」で、S114で求めたゾーン位置が「第1文字目～第500文字目」であるので、キーワード検索手段6は、文書識別子ID1とそれが第1番目の検索項目についてのものであることを示す情報とをペアにしてキーワード検索結果格納手段8に格納することになる。

【0049】その後、キーワード検索手段6は、第2番目の検索項目「要約=インデックス」に注目し(S111)、前述したと同様の処理を行う(S113～S115)。第2番目の検索項目の場合、キーワード「インデックス」は、文書識別子ID1のゾーン位置変換文書21'の第1001文字目～第1006文字目に現れ、ゾーン名「要約」のゾーンのゾーン位置は、「第501文字目～第2000文字目」であるので、キーワード検索手段6は、文書識別子ID1とそれが第2番目の検索項目についてのものであることを示す情報とをペアにして

キーワード検索結果格納手段8に格納する。

【0050】そして、検索条件入力手段7から渡された全ての検索項目について上述した処理を行うと(S112がNO)、キーワード検索手段6は、論理条件解析手段9に対して終了通知を送る(S116)。

【0051】論理条件解析手段9は、キーワード検索手段6から終了通知が送られてくると、図10の流れ図に示すように、キーワード検索結果格納手段8に格納されている各検索項目についての検索結果と、検索条件入力手段7から渡された検索条件式中の論理演算記号とに基づいて、検索条件式を満足させる構造化文書の文書識別子を求め、それを検索結果出力手段10に渡す(S101、S102)。

【0052】この例の場合、キーワード検索結果格納手段8には、第1番目、第2番目の検索項目の検索結果としてそれぞれ文書識別子「ID1」、「ID1」が格納され、検索条件式中の第1番目の検索項目と第2番目の検索項目とを結合する論理演算式が「AND」であることから、両方の検索結果中に存在する文書識別子「ID1」を検索結果出力手段10に渡すことになる。

【0053】検索結果出力手段10は、文書識別子「ID1」が渡されると、文書格納手段1から文書識別子が「ID1」の構造化文書21を読み込み、プリンタ、CRT等の出力装置(図示せず)に出力する。

【0054】図13は、図1に示した構造化文書検索装置のハードウェア構成の一例を示したブロック図であり、コンピュータ131と、記録媒体132と、記憶装置133とから構成されている。記録媒体132は、磁気ディスク、半導体メモリ、その他の記録媒体であり、コンピュータ131を構造化文書検索装置として機能さ

せるためのプログラムが記録されている。

【0055】記録媒体132に記録されているプログラムは、コンピュータ131によって読み込まれ、コンピュータ131の動作を制御することにより、コンピュータ131上に図1に示した文書内位置変換手段2、インデックス作成手段4、キーワード検索手段6、検索条件入力手段7、論理条件解析手段9、検索結果出力手段10を実現する。尚、文書格納手段1、ゾーン情報テーブル3、インデックス5、キーワード検索結果格納手段8は、記憶装置133上に構成される。

【0056】

【発明の効果】以上説明したように、本発明の構造化文書検索装置によれば、従来の全文インデックスを利用してゾーン検索を行う従来の技術に比較してインデックスサイズを小さくすることができ、且つ検索速度を高速化することができる。その理由は、ゾーンの位置情報が全ての構造化文書で共通になるような形でインデックスを作成するため、検索時に、非常に小規模なゾーン情報テーブルを参照するだけでゾーン検索を行うことができるからである。

【図面の簡単な説明】

【図1】本発明の実施例のブロック図である。

【図2】文書格納手段1の内容例を示す図である。

【図3】ゾーン情報テーブル3の内容例を示す図である。

【図4】インデックス5の内容例を示す図である。

【図5】検索条件式の一例を示す図である。

【図6】文書内位置変換手段2の処理例を示す流れ図である。

【図7】インデックス作成手段4の処理例を示す流れ図である。

【図8】検索条件入力手段7の処理例を示す流れ図である。

【図9】検索条件入力手段7から検索条件式が渡されたときの論理条件解析手段9の処理例を示す流れ図である。

【図10】キーワード検索手段6から終了通知が送られてきたときの論理条件解析手段9の処理例を示す流れ図である。

【図11】キーワード検索手段6の処理例を示す流れ図である。

【図12】文書内位置変換手段2で作成されたゾーン位置変換文書の一例を示す図である。

【図13】図1に示した構造化文書検索装置のハードウェア構成の一例を示すブロック図である。

【符号の説明】

1…文書格納手段

2…文書内位置変換手段

3…ゾーン情報テーブル

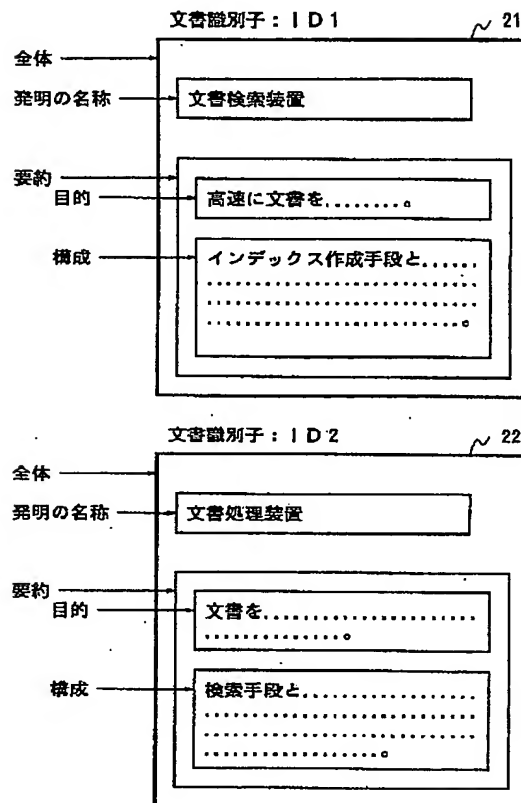
4…インデックス作成手段



14

- 2 1', 2 2' …ゾーン位置変換文書  
5 1…キー情報部  
5 2…位置情報部  
1 3 1…コンピュータ  
1 3 2…記録媒体  
1 3 3…記憶装置

【图 2】



【図 3】

3 ゾーン情報テーブル	
ゾーン名	ゾーン位置情報
全体	1文字目 ～ 2000文字目
発明の名称	1文字目 ～ 500文字目
要約	501文字目 ～ 2000文字目
目的	501文字目 ～ 1000文字目
構成	1001文字目 ～ 2000文字目

【図 5】

ゾーン名      キーワード                  ゾーン名      キーワード

|                    |                    |                    |

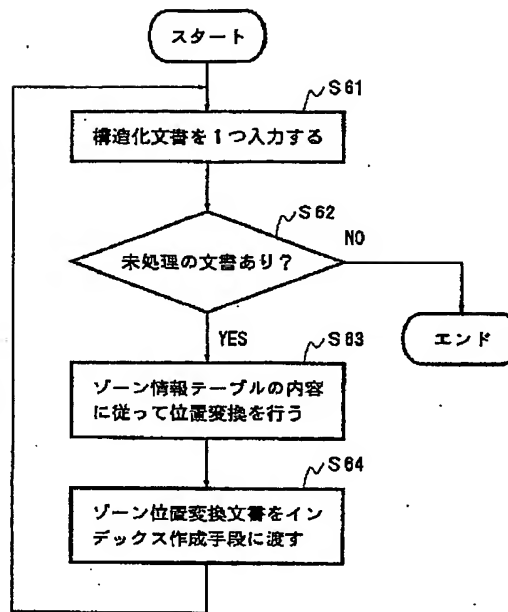
(発明の名称=検索)   AND   (要約=インデックス)

【図4】

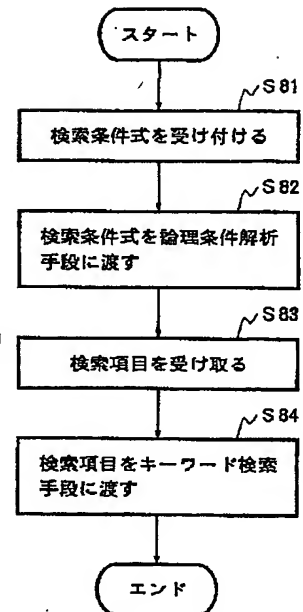
5 インデックス

キー情報	位置情報 (文書識別-文書内位置)
文	1-1, 1-504, 2-1, 2-501
書	1-2, 1-505, 2-2, 2-502
検	1-3, 2-1001
索	1-4, 2-1002
装	1-5, 2-5
置	1-6, 2-6
高	1-501
速	1-502
⋮	⋮
イ	1-1001
ン	1-1002
デ	1-1003
ッ	1-1004
ク	1-1005
ス	1-1006
⋮	⋮

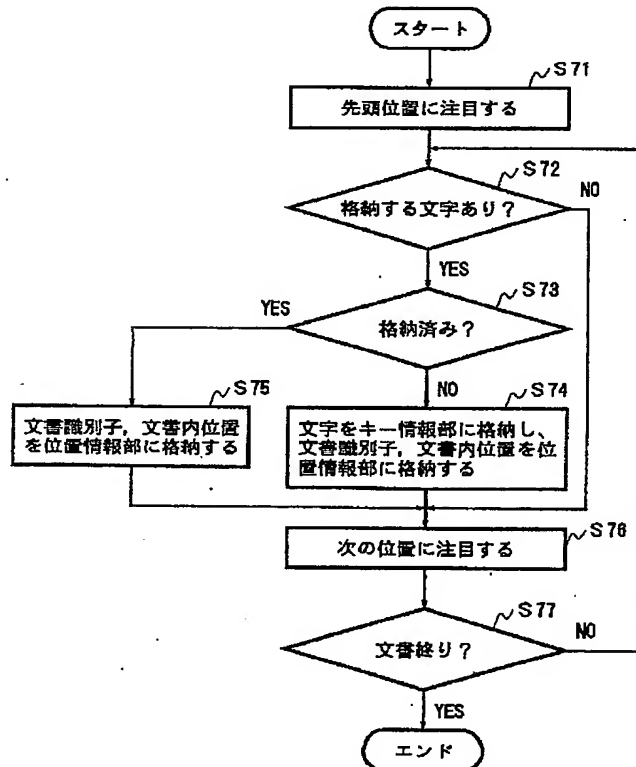
【図6】



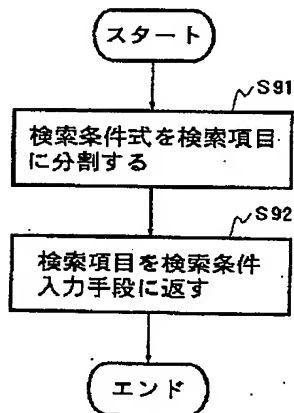
【図8】



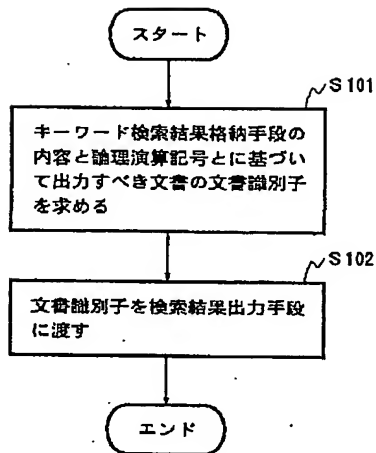
【図7】



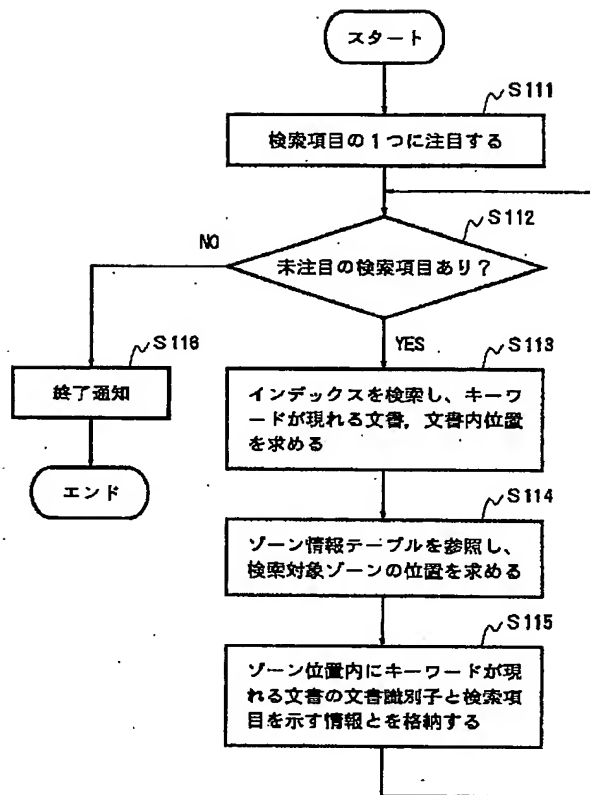
【図9】



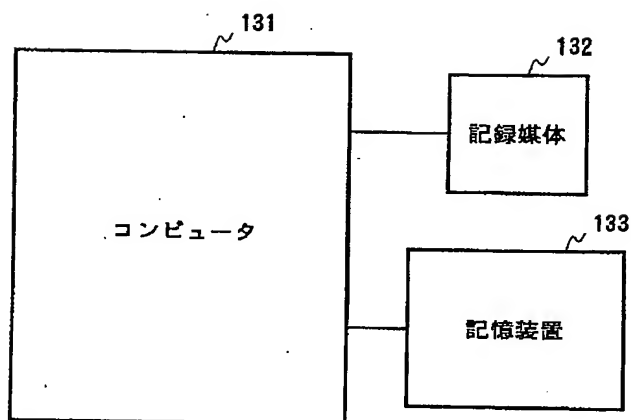
【図 10】



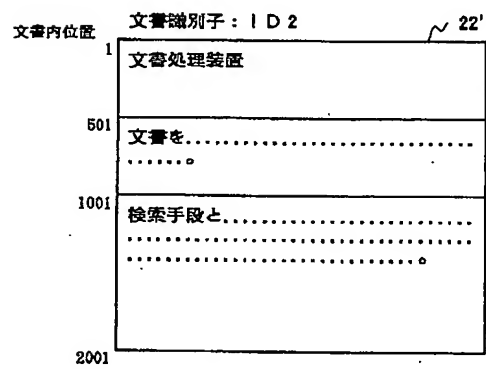
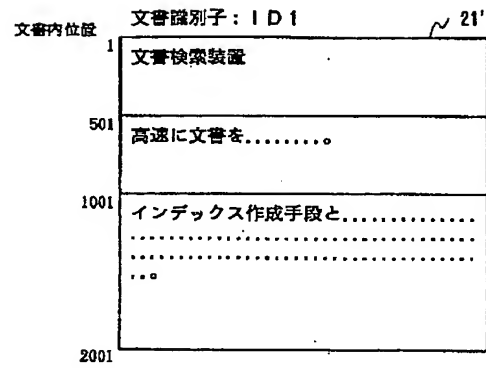
【図 11】



【図 13】



【図 1 2】



**This Page is Inserted by IFW Indexing and Scanning  
Operations and is not part of the Official Record**

**BEST AVAILABLE IMAGES**

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☐ BLACK BORDERS
- ☐ IMAGE CUT OFF AT TOP, BOTTOM OR SIDES
- ☒ FADED TEXT OR DRAWING
- ☒ BLURRED OR ILLEGIBLE TEXT OR DRAWING
- ☐ SKEWED/SLANTED IMAGES
- ☐ COLOR OR BLACK AND WHITE PHOTOGRAPHS
- ☐ GRAY SCALE DOCUMENTS
- ☐ LINES OR MARKS ON ORIGINAL DOCUMENT
- ☒ REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY
- ☐ OTHER: \_\_\_\_\_

**IMAGES ARE BEST AVAILABLE COPY.**

**As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.**